

Resurse lingvistice pentru un sistem de întrebare-răspuns pentru limba română

Verginica Barbu Mititelu, Alexandru Ceaușu, Radu Ion,
Elena Irimia, Dan Ștefănescu, Dan Tufiș

Institutul de Cercetări pentru Inteligență Artificială, Academia Română
Calea 13 Septembrie, nr. 13, 050711, București
E-mail: {vergi, aceausu, radu, elena, danstef, tufis}@racai.ro

Rezumat. În acest articol descriem câteva resurse lingvistice (un lexicon, o morfologie paradigmatică, două tezaure lingvistice – wordnetul românesc și Eurovoc – și un corpus multilingv paralel) din perspectiva utilității lor preponderent în sistemele de întrebare-răspuns. Prezentăm aici etapele procesului de găsire automată a unui răspuns la o întrebare formulată în limbaj natural de utilizator. Acolo unde se impune, este înfățișat modul în care resursele lingvistice contribuie la rezolvarea diverselor probleme. Lexiconul funcționează ca un corector ortografic pentru întrebarea introdusă de utilizator. Morfologia paradigmatică este utilă în procesul de lematizare a întrebării și a corpusului. Wordnetul românesc își dovedește eficacitatea în expandarea frazei de interogare, în identificarea lanțurilor lexicale între sensurile cuvintelor și în găsirea răspunsurilor într-un sistem mono-și multilingv. Tezaurul Eurovoc este folosit la segmentarea și lematizarea întrebării introduse de utilizator și a corpusului multilingv paralel în care se caută răspunsul la această întrebare. Arhitectura sistemului de întrebare-răspuns descris aici nu este dependentă de limbă; resursele lingvistice, însă, sunt, în mod intrinsec, specifice unei limbi (vezi lexiconul și morfologia paradigmatică); excepție fac cele a căror organizare sau structură permite o perspectivă multilingvă (ex. tezaurele și corpusul), în cazul nostru ele fiind aliniate între ele.

Cuvinte cheie: sistem de întrebare-răspuns, resurse lingvistice, tezaur, ontologie lexicală, corpus.

1. Introducere

Interesul pentru lingvistica computațională și pentru finanțarea proiectelor din acest domeniu de cercetare este susținut de aplicațiile pe care le dezvoltă, cu impact pentru societate în ansamblul ei. Lumea de astăzi are, prin internet, acces la o cantitate de informații extrem de vastă. Sunt necesare instrumente care să permită găsirea datelor petinente în acest volum uriaș de cunoștințe. Comunitatea lingvisticii computaționale ține

pasul cu aceste nevoi și încearcă să ofere sisteme informatice menite să ajute utilizatorul să se descurce în noianul de documente disponibile.

Pentru informare, utilizatorul recurge frecvent la formularea unei întrebări. Pornind de la aceasta, diverse aplicații oferă variate soluții: documentele care, se presupune, includ răspunsul la întrebare sau chiar fragmentele considerate adecvate din documentele estimate a fi cele mai relevante.

Regăsirea informației (engl. *information retrieval*) este știința aspectelor teoretice și practice ale căutării informațiilor de interes în colecții mari de date, ale creării și menținerii acestor colecții (Jackson și Moulinier, 2002 : 26).

Extragerea informației (engl. *information extraction*) este un tip de regăsire a informației, care-și propune să identifice, în mod automat, anumite tipuri de entități, relații, evenimente în texte. Se deosebește de regăsirea informației prin faptul că nu își propune găsirea de documente, ci de informație relevantă în interiorul documentelor (Jackson și Moulinier, 2002 :75).

Sistemele de întrebare-răspuns (engl. *question answering*) găsesc, în mod automat, și redau răspunsuri la întrebări formulate în limbaj natural (Harabagiu și Moldovan, 2003 : 561). De la utilizatori ai Internetului dornici să-și îmbogățească cunoștințele generale, la analiști specialiști într-un domeniu, care vor să afle o informație specifică, pentru toți astfel de sisteme devin necesare și de neînlocuit.

2. Utilitatea resurselor lingvistice în aplicațiile lingvisticii computaționale

Aplicațiile de regăsire și extragere a informației și cele de întrebare-răspuns au la dispoziție doar cuvintele introduse de utilizator în întrebare. În primul rând, trebuie verificată ortografia acestora, iar pentru aceasta este util un *lexicon*. În al doilea rând, cuvintele introduse de utilizator pot să nu se regăsească în documente exact în forma din întrebare. În astfel de cazuri își dovedește utilitatea un *lematizor*. În al treilea rând, cuvintele folosite de utilizator nu sunt neapărat cele mai relevante pentru găsirea răspunsului. De aceea, este necesară expandarea frazei de interogare, prin luarea în considerare a cuvintelor aflate în relații semantice sugestive (sinonimie, hiponimie, hiperonimie) cu cele care compun această frază de interogare.

Pentru aceasta este nevoie de un *tezaur*, în care sunt înregistrate cuvinte cu sensurile lor și cu relațiile lexico-semantice existente între ele. În al patrulea, dar în nici un caz ultimul rând, este necesară o colecție de documente în care să fie căutat răspunsul. Această colecție poate fi un *corpus* sau internetul.

În acest articol, prezentăm, din perspectiva utilității (preponderent) în sistemele de întrebare-răspuns, cinci resurse lingvistice disponibile la Institutul de Cercetări pentru Inteligență Artificială (ICIA) al Academiei Române: un lexicon, o morfologie paradigmatică pentru limba română, un tezaur lingvistic general (ontologie lexicală), un tezaur specializat și un corpus de mari dimensiuni.

2.1 Lexicon

Motoarele de căutare iau cuvintele din întrebarea introdusă de utilizator și întorc rezultate ordonate pentru aceasta. În cadrul acestui proces, un factor important care duce la rezultate nesatisfăcătoare este incorectitudinea scrierii cuvintelor de către utilizator. Se pare că 10-15% dintre frazele de căutare conțin astfel de greșeli (Cucerzan și Brill, 2004). De multe ori este vorba despre erori dactilografice, însă cele mai dăunătoare sunt cele rezultate din confuzia unor termeni: paronimia (i.e. relația dintre două cuvinte asemănătoare din punctul de vedere al formei, dar deosebite ca sens), de pildă, poate crea confuzii numeroase: *revela* și *releva*. În funcție de colecția de documente în care se face căutarea, se preferă diverse modalități de corectare a acestor erori. În cazul nostru, al căutării într-un corpus¹ își poate dovedi utilitatea un lexicon cuprinzător. Fiecare element din întrebare (supusă în prealabil unei operații de segmentare la nivel de cuvânt) este verificat din punct de vedere ortografic prin căutarea formei în lexicon. Dacă aceasta nu este găsită, se aplică o euristică de găsire a celei mai apropiate forme pe baza unor principii simple: erori de dactilografiere

¹ Pentru căutarea răspunsurilor pe web, se impun alte modalități de corectare a erorilor de scriere, determinate de vastitatea materialului lexical ce trebuie acoperit și de raritatea utilizării anumitor cuvinte (vezi Chen et al. 2007).

(inversarea unor litere apropiate pe tastatură), litere alăturate în cuvânt inversate, litere duble.

Lexiconul românesc a fost creat în cadrul proiectului MULTEXT-EAST și conține forme flexionare ale cuvintelor din limba română, împreună cu descrierea lor morfosintactică și cu lema (i.e. forma de dicționar, Tușiș et al., 1997, vezi Figura 1). După încheierea proiectului, în ICIA s-a continuat îmbogățirea acestui lexicon, iar în momentul de față, el conține aproape 1.223.000 de intrări (i.e. forme flexionare), față de 450.000 câte conținea la finalizarea MULTEXT-EAST (cf. Tușiș, 2002).

Formă ocurentă	Lemă	MSD	CTAG
aramă	aramă	Ncfsrn	NSRN
arama	aramă	Ncfsry	NSRY
arama	aramă	Ncfp-n	NPN
arama	aramă	Ncfson	NSON
arama	aramă	Ncfsoy	NSOY
arama	aramă	Ncfpry	NPRY
arama	aramă	Ncfpoy	NPOY

Figura 1. Formele flexionare ale substantivului *aramă* în lexicon.

O situație aparte care impune corectură grafică o reprezintă introducerea de către utilizator a cuvintelor fără diacritice sau numai parțial cu diacritice. Cum corpusul în care motorul nostru face căutarea (vezi *infra* 2.4) conține diacritice, se impune introducerea acestora și în întrebare. În cadrul ICIA a fost implementat sistemul DIAC (Tușiș și Chițu, 1999) extins și mult îmbunătățit de sistemul DIAC⁺, care are o precizie de 97.75% în reconstituirea diacriticelor în texte neprocesate (Tușiș și Ceaușu, 2008).

Principala problemă ridicată de utilizarea lexiconului în corectarea grafiei o reprezintă limitarea cantitativă a acestuia. În consecință, inexistența unui cuvânt din întrebare în lexicon nu trebuie interpretată neapărat ca greșeală de scriere.

2.2 Morfologie paradigmatică pentru limba română

Odată stabilită corectitudinea grafică a întrebării utilizatorului urmează adnotarea de tip stratificat (Tușiș, 1999) cu etichete morfosintactice (vezi Ion, 2007 pentru descrierea completă a acestui proces) și lematizarea ei,

adică identificarea formei de bază (lema) a fiecărui cuvânt, cu ajutorul lexiconului. Dată fiind o formă ocurentă și descrierea sa morfosintactică, se regăsește ușor lema formei respective în lexiconul de forme flexionate. Pentru cazurile neacoperite de acesta, la ICIA a fost dezvoltat un lematizor (Ion, 2007: 22-26) care folosește un set de reguli (specific fiecărei categorii gramaticale flexionare) induse automat din lexicon, care generează leme candidat pentru cuvântul necunoscut, și apoi modele Markov (antrenate pe leme din lexicon) pentru a ordona candidații. Candidatul cu cea mai mare probabilitate câștigă. Procedura funcționează foarte bine, cele mai multe dintre erori sunt în cazul cuvintelor necunoscute care aparțin paradigmelor flexionare neregulate sau atunci când adnotarea morfosintactică a formei ocurență a fost greșită. În ansamblu, ținând cont de vasta acoperire a lexiconului și de rata mică de eroare a lematizorului statistic, probabilitatea unei erori de lematizare este neglijabilă.

În vederea căutării răspunsului la întrebarea utilizatorului, se impune și preprocesarea corpusului. El este supus, pe rând, operațiilor de segmentare, adnotare cu etichete morfosintactice și lematizare. Această ultimă etapă se face tot cu ajutorul lematizorului descris mai sus. Pentru formele inexistente în lexicon, însă, mai dispunem de o modalitate de identificare a lemei, cu ajutorul morfologiei paradigmatică pentru limba română (ROPMORPH, Tufiș, 1989, Irimia, 2009a), care conține un inventar complet al terminațiilor din limba română pentru verbe, substantive și adjective, fiecare terminație fiind asociată informației morfologice caracteristice unei anumite forme (vezi Figura 2).

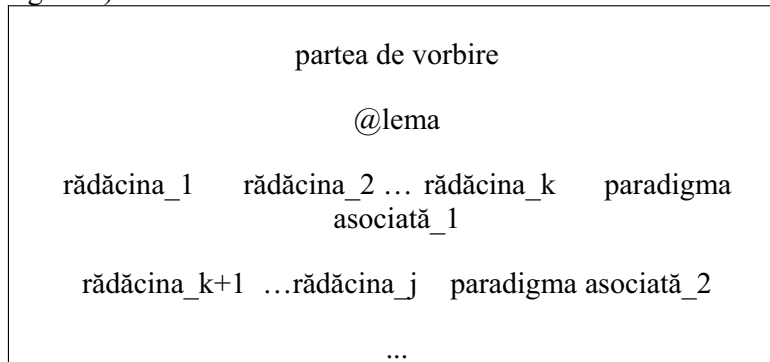


Figura 2. Structura unei intrări în ROPMORPH.

Astfel, pentru fiecare formă flexionară, poate fi obținută o listă de rădăcini prin eliminarea terminațiilor specifice fiecărei paradigme aplicabile formei. În cazul în care lista conține un singur element, cu ajutorul generatorului morfologic (ROG) (Irimia, 2009b) se poate, foarte ușor, genera lema corespunzătoare formei flexionare de la care am pornit. În cazul în care lista conține două sau mai multe rădăcini, Irimia (2009a) descrie un algoritm prin care ROG generează toate formele posibile din toate paradigmele asociate acestor rădăcini, apoi apelează o serie de funcții dintr-o bibliotecă dezvoltată la ICIA, care invocă motorul de căutare GoogleTM și extrag informații cantitative, pentru fiecare cuvânt din fiecare paradigmă, referitoare la numărul de apariții pe web. Sunt eliminate toate formele care apar de mai puțin de 10 ori (se presupune că ele se datorează unor greșeli de scriere), iar fiecare paradigmă este ierarhizată în funcție de suma ocurențelor pe web ale formelor sale. Este aleasă paradigma cu cel mai mare scor, iar apoi ROG generează lema.

2.3. Resurse lexico-semantice

Tezaurul este un tip de dicționar în care cuvintele cu sensuri asemănătoare sunt grupate la un loc (cf. Online Cambridge Dictionary). Prima lucrare de acest fel, din epoca modernă, a fost creată de Peter Mark Roget în 1805 și cuprindea 15.000 de cuvinte. De atunci a cunoscut adăugiri succesive, la început de către urmași ai lui Roget, apoi de către alți specialiști, rezultatul fiind câte o ediție nouă a acestei lucrări, ultima dintre ele (din 1987) cuprinzând peste 1.250.000 de cuvinte.

Organizarea materialului lingvistic se face în clase, fiecare având o ierarhie internă, asemănătoare unui arbore cu multiple ramificații pentru grupuri de sensuri. Scopul unei asemenea lucrări este să faciliteze găsirea cuvintelor înrudite semantic (sinonime, antonime, hiperonime, hiponime), ceea ce în dicționarele explicative este mai greu sau chiar imposibil, și să ajute la alegerea celui mai potrivit cuvânt pentru exprimarea sensului dorit de utilizator.

În prelucrarea limbajului natural se folosesc în ultimii ani, din ce în ce mai mult, ontologiile lexicale. Ca și în tezaurele binecunoscute în regăsirea de documente, în ontologiile lexicale sunt stocate numeroase cuvinte (în forma lor lematizată), împreună cu relațiile semantice sau lexicale dintre ele. Cuvintele sunt abstractizate ca unități semantice pentru exprimarea

conceptelor. Între ele se stabilesc relații ierarhice (pentru indicarea termenilor mai generali și a celor mai specifici), de echivalență (stabilite între sinonime) sau asociative (relații care nu sunt nici ierarhice, nici de echivalență, dar se stabilesc între cuvinte, iar utilizatorul le acceptă ca răspuns la întrebarea sa). Ceea ce diferențiază o ontologie lexicală de un tezaur este nivelul axiomatic al descrierii conceptelor denotate de cuvintele reprezentate. Pe lângă definiție, relațiile de tip ierarhic (hiperonim, hiponim), meronimic sau de altă natură, o ontologie lexicală pune în evidență o legătură de tip ontologic, la un concept interlingual caracterizat de o serie de definiții și axiome. De exemplu, pentru cuvântul „container”, descrierea în ontologia SUMO/MILO (Niles și Pease, 2001) este prezentată în Figura 3.



Figura 3: Descrierea axiomatică a conceptului interlingual „container”

2.3.1 RO-Wordnet, ontologie lingvistică pentru limba română

Lexicul este, fără îndoială, cea mai importantă resursă lingvistică a unei limbi. Marea majoritate a cercetărilor actuale, atât în lingvistica formală, cât mai ales în tehnologia limbajului, plasează componenta lexicală în centrul

modelelor de limbă, sub influența a ceea ce a fost numită abordarea *lexicalizată* sau *lexicalistă* a studiului limbii. Studiul computațional al dicționarelor electronice, natura informației ce trebuie inclusă în ele și tipul de prelucrări pe care le poate facilita o anumită structurare a unui mare volum lexical a fost, fără îndoială, fundamental influențat de proiectul WordNet, lansat în 1985 la Universitatea din Princeton sub conducerea reputatului psiholingvist George Miller (Fellbaum, 1998). WordNet, resursă publică, în varianta actuală, este o uriașă rețea semantică lexicală, în care peste 117.000 de *înțelesuri* lexicalizate în limba engleză prin mai mult de 155.000 de cuvinte sunt asociate între ele prin relații semantice (Fellbaum, 1998). Fondul lexical este distribuit în patru rețele semantice corespunzând categoriilor gramaticale deschise: substantive, verbe, adjective și adverbe. Noțiunea de *înțeles* (engl. *meaning*) este echivalată în WordNet cu cea de *concept* și este reprezentată printr-o serie sinonimică în care fiecărui cuvânt al seriei îi este asociat un număr ce identifică sensul în care cuvântul respectiv are înțelesul asociat conceptului. Seria sinonimică ce identifică un înțeles se numește *sinset* (engl. *synonymy set*). Relațiile existente între *sinseturi* sunt de diferite tipuri, depinzând de categoria gramaticală a cuvintelor ce alcătuiesc un anumit sinset (antonimie/sinonimie, hiponimie/hiperonimie, holonimie/meronimie, troponimie etc.) și având justificare psiholingvistică.

Wordnetul este un dicționar semantic care rafinează structurile întâlnite în tezaure, prin considerarea aspectelor lexicale ale materialului lingvistic conținut: de exemplu, rețeaua semantică a substantivelor cuprinde 25 de structuri arborescente (relațiile organizatoare fiind hiponimia și meronimia), iar cea a verbelor 15 astfel de structuri (organizate în funcție de hiponimie și troponimie), asemănătoare claselor din *Tezaurul* lui Roget. Adjectivele și adverbele nu se pretează, însă, la astfel de organizări ierarhice: în cazul adjectivelor se preferă gruparea lor în funcție de relația de antonimie, iar adverbele (cu excepția cazurilor rare de antonimie) nu au nicio organizare.

Influența proiectului WordNet a fost enormă în domeniul tehnologiei limbajului (exprimată poate și prin faptul că acum, în limbajul tehnic cel puțin, cuvintele „wordnet” și „synset” au devenit substantive comune, importate prin calchiere în mai toate limbile), iar beneficiile acestui concept atât de evidente încât Comisia Europeană a decis finanțarea a două proiecte de mare anvergură: EuroWordNet (Blokma et al., 1996) și BalkaNet (Stamou et al., 2001).

În cadrul BalkaNet a început dezvoltarea (Tufiş şi Barbu, 2004) ontologiei lingvistice pentru limba română (Ro-WN), păstrându-se organizarea materialului lingvistic din wordnetul englezesc (Tufiş, 2004). Consolidarea acestei resurse lingvistice a continuat în cadrul programului naţional „Cercetare de Excelenţă” finanţat de Ministerul Educaţiei şi Cercetării prin proiectul ROTEL (<http://www.ici.ro/~badea/rotel/index.htm>) şi ca temă de plan a Institutului. În momentul scrierii acestei lucrări, Ro-WN conţine peste 55.000 de sinseturi şi peste 47.000 de literalii unici: 40.495 de sinseturi substantivale (adică 73%), 9.940 de sinseturi verbale (18%), 3831 de sinseturi adjectivale (7%) şi 834 de sinseturi adverbiale (2%). Toate sinseturile conţinute în Ro-WN sunt aliniată cu echivalenţii de traducere din limba engleză din Princeton WordNet (versiunea 2.0) şi, în plus, au legături ontologice la conceptele din ontologia de nivel superior şi mediu SUMO/MILO (vezi Figura 3). Sinseturile ontologiei lexicale pentru limba română beneficiază şi de marcaje de tipul clasificării zecimale universale (Bentivogli et al., 2004) şi pentru specificarea triadei de subiectivitate (Esuli şi Sebastiani, 2006) <obiectiv, pozitiv, negativ> (a se vedea pentru detalii Tufiş, 2008).

În cadrul unui sistem de întrebare-răspuns, o ontologie lexicală precum wordnetul îşi dovedeşte utilitatea în mai multe etape ale lucrului: expandarea frazei de interogare, identificarea lanţurilor lexicale între sensurile cuvintelor şi găsirea răspunsurilor într-un sistem monolingv sau multilingv.

2.3.1.1 Expandarea frazei de interogare

După prelucrare, întrebarea este transformată în frază de interogare în metalimbajul sistemului de căutare. De exemplu, întrebarea utilizatorului

„Care este scopul modificării deciziei 77/270/Euratom în ceea ce priveşte energia nucleară?”

devine

„label:REASON +("scop modificare decizie 77/270"^4 (scop AND modificare AND decizie AND 77/270) Euratom energie_nucleară "scop modificare decizie"^3 (scop AND modificare AND decizie) "modificare decizie 77/270"^3 (modificare AND decizie AND 77/270) "scop modificare"^2 (scop AND modificare) "modificare decizie"^2 (modificare AND decizie) "decizie 77/270"^2 (decizie AND 77/270) scop modificare

decizie 77/270 energie_nucleară)". Sistemul caută un paragraf etichetat ca fiind de tip „scop” (label:REASON), în care să existe „scop”, „modificare”, „decizie”, „77/270” etc. Notațiile de tip \wedge^2 reprezintă un coeficient de creștere a relevanței (engl. *boosting*) în căutare a termenului căruia i se aplică².

Întrebarea introdusă de utilizator poate prezenta deficiențe și în ceea ce privește proprietatea termenilor utilizați (i.e. folosirea acelor cuvinte care să reflecte cel mai bine sensul). Din acest motiv, pentru găsirea răspunsului, este utilă expandarea frazei de interogare, prin includerea în căutare a cuvintelor sinonime și chiar hiperonime, hiponime sau troponime cu cele introduse de utilizator. Organizarea ontologiei lexicale slujește foarte bine acestui scop: wordnetul conține sinseturi (grupuri de sinonime) ierarhizate conform relațiilor de hiponimie și troponimie. Folosind un sistem de dezambiguizare automată a sensurilor cuvintelor din fraza de interogare (Ion și Tușiș, 2007), se caută în Ro-WN cuvintele aflate în relații semantice relevante cu acestea (în special sinonimie, hipo- și hiperonimie, troponimie) și sunt luate în considerare în căutarea răspunsului în corpus.

2.3.1.2 Identificarea lanțurilor lexicale între sensurile cuvintelor

Pentru a putea ordona mulțimea răspunsurilor găsite de sistem pentru o frază de interogare, este nevoie să se stabilească un scor de apropiere semantică între cuvintele introduse de utilizator și cele materializate în text (acestea nefiind întotdeauna identice, motiv pentru care se efectuează și expandarea descrisă în 2.3.1.1). Calcularea acestui scor de apropiere semantică se face pe baza lungimii lanțurilor lexicale dintre cuvinte, folosind ontologia wordnet pentru limba română.

Un lanț lexical reprezintă o mulțime de cuvinte înrudite semantic, specifice unui domeniu de discurs. Deosebirea față de cuvintele care formează un câmp semantic o reprezintă faptul că un lanț lexical trebuie interpretat ca o cale în ierarhia Ro-WN, cale ce include, cel mai adesea, și relații intercategoriale ce dau seama de dependențele sintagmatice ce pot

² Pentru descrierea componentelor, a modului de funcționare și a ordonării rezultatelor sistemului de întrebare-răspuns dezvoltat în cadrul ICIA, vezi Ion et al. (2008).

exista între cuvintele care definesc capetele unui lanț lexical (vezi Ion et al., 2008) pentru descrierea acestei metode de ordonare a răspunsurilor). Diverse relații participă în mod diferit la stabilirea scorului de asemănare între două cuvinte: hiponimia și relațiile derivatice (intercategoriale) au cea mai importantă contribuție.

În exemplul de mai jos redăm trei lanțuri lexicale găsite între *fiu* și *tată*, precedate de scorurile de asemănare (*near-antonym* notează relația de antonimie, *gloss* indică ocurența în definiție, *hypernym* notează relația de hiperonimie, iar *ss* relația de sinonimie, deci coocurența în același sinset):

0.333

```
fiu(?) near_antonym fiică(1) <=> fiică(1) near_antonym copil(?)
gloss tată(?)
```

0.666

```
fiu(?) hypernym băiat(1) <=> băiat(1) hypernym copil(?) gloss
tată(?)
```

0.733

```
fiu(?) ss băiat(5) <=> băiat(5) ss copil(?) gloss tată(?)
```

2.3.1.3 Găsirea răspunsurilor într-un sistem de întrebare-răspuns multilingv

Imaginând un scenariu în care utilizatorul vrea să afle răspunsul într-o limbă străină (limba țintă pe care nu o stăpânește foarte bine) la o întrebare pe care preferă să o formuleze în limba maternă (limba sursă), se poate dovedi utilitatea resurselor lingvistice multilingve aliniat. Este vorba despre tezaure lingvistice sau ontologii lexicale multilingve aliniat între ele și despre corpusuri aliniat la nivel de paragraf sau chiar de propoziție (vezi *infra* 2.4).

Procedura de lucru este identică, până într-un punct, cu cea a sistemului de întrebare-răspuns monolingv. Întrebarea este prelucrată, fraza de interogare este expandată și se caută răspunsul în corpusul în limba sursă. Cum Acquis-ul comunitar (vezi *infra* 2.4 pentru detalii) în limba română este aliniat la nivel de paragraf cu Acquis-ul comunitar al tuturor țărilor membre ale Uniunii Europene, se poate identifica foarte rapid paragraful corespunzător din limba țintă (Ceașu et al., 2006).

O altă modalitate de găsire a răspunsului într-o limbă țintă este de a traduce automat întrebarea utilizatorului și de o a supune prelucrării, expandării și apoi de a căuta răspunsul în corpusul din acea limbă. Această

metodă, însă, va da rezultate a căror relevanță este diminuată de calitatea sistemului de traducere automată.

În sfârșit, o a treia variantă de rezolvare a problemei de interogare într-o limbă și de căutare a răspunsului în documente elaborate în altă limbă este de a traduce automat doar grupurile sintactice (nu toată întrebarea) considerate relevante în construcția frazei de interogare. Acest lucru revine la a traduce fraza de interogare construită pentru limba sursă. La ICIA au fost deja experimente în acest sens, calitatea traducerilor grupurilor nominale din engleză în română și invers fiind extrem de bună. Una din direcțiile de dezvoltare a sistemului de întrebare-răspuns al ICIA, pe lângă transpunerea sa într-un serviciu web public, este încorporarea acestei abordări pentru tratarea situațiilor interlinguale de întrebare-răspuns.

2.3.2 Eurovoc

Eurovoc este un tezaur multilingv (Steinberger et al., 2002), folosit la indexarea Acquis-ului comunitar. El conține termeni din domeniile de activitate din țările membre UE și relații între aceștia. Eurovoc versiunea 4.3 există în 23 de limbi, acoperă 21 de domenii și este organizat în 127 de microtezaure.

Termenii sunt de două feluri: descriptori (termenii folosiți în indexare) și non-descriptori (cuvinte din limba comună echivalente cu descriptorii, dar nefolosite în indexare; sunt folosite la expandarea frazei de interogare). Relațiile semantice sunt de cinci tipuri: cele care definesc descriptorii, cele care indică apartenența acestora la unul sau mai multe microtezaure, relații de echivalență (stabilite între descriptorii și non-descriptorii), relații ierarhice (stabilite între descriptorii) și asociative (reprezentând modalitățile în care descriptorii se pot asocia: cauză, efect, locație etc.)

Fiecare document din Acquis-ul comunitar este indexat cu termenii Eurovoc pe care îi conține. Acest lucru facilitează găsirea răspunsului, prin limitarea căutării la acele documente care conțin descriptorii din întrebarea utilizatorului sau non-descriptorii aflați în diverse relații cu aceștia.

Sistemul de întrebare-răspuns dezvoltat la ICIA (Ion et al., 2008) face apel la Eurovoc în segmentarea și lematizarea întrebării utilizatorului și a

corpusului. Termenii din Eurovoc sunt considerați unități minimale în etapa de segmentare, iar la lematizare li se atașează lema din Eurovoc³.

2.4 Corpusul paralel Acquis comunitar

Nevoia de corpusuri paralele pentru aplicațiile de procesare a limbajului natural a cunoscut o accentuare pe parcursul ultimilor ani. Corpusurile paralele sunt folosite în aplicațiile de traducere automată sau categorizare multilinguală, pentru a produce resurse lexicale sau semantice multilinguale, cum sunt dicționarele sau ontologiile lexicale, pentru a testa consistența procesului de traducere etc.

Pentru aplicațiile de procesare a limbii române în context multilingual cel mai important corpus este JRC-Acquis (Steinberger et al., 2006). Acesta este, în prezent, cel mai mare corpus multilingual disponibil, conținând 22 de limbi. Corpusul este disponibil în format XML conform specificațiilor TEI (Text Encoding Initiative). De asemenea, conține și alinierea celor peste 231 de perechi de limbi conținute de JRC-Acquis. Corpusul crește pe măsura traducerii legislației europene și în limbile noilor candidați.

Acquis-ul comunitar este corpul comun de convenții, legi și obligativități care leagă toate statele membre ale Comunității Europene. Conține principii și obiective politice ale diverselor tratate semnate în cadrul Uniunii Europene, legislație UE, declarații și rezoluții, acorduri internaționale și obiective comune. Pe lângă cele 22 de limbi ale Comunității Europene, Acquis-ul este tradus și în limbile croată și turcă. Datorită efortului depus la ICIA pentru colectarea și adnotarea documentelor Acquis-ului românesc, limba română a fost prezentă în pachetul de distribuție JRC-Acquis încă de la prima versiune a acestuia.

Documentele disponibile în limba română sunt 19.211. Numărul de documente comune în perechea de limbi engleză-română este de 11.469. Documentele comune constituie un important corpus paralel conținând 59.986.838 de cuvinte.

³ În alegerea sinseturilor de implementat în Ro-WN atenție deosebită a fost acordată acoperirii integrale a tezaurului Eurovoc.

Fișierele românești și bulgărești nu sunt disponibile în același format (HTML) ca documentele celorlalte limbi din Acquis, neputând fi procesate de aceleași instrumente de conversie HTML-TEI. Fișierele în format HTML disponibile pentru celelalte limbi ale JRC-Acquis conțin și informații cu privire la structura documentului, cum ar fi secțiunile de anexe și semnături, secțiunile cu textul și titlul documentului etc. Această structură nu se regăsește în formatul Microsoft Word, în care sunt disponibile documentele românești ale JRC-Acquis.

Pentru a constitui colecția de documente în limba română, fișierele au fost descărcate de pe site-ul „CCVista Translation Database” (<http://ccvista.taix.be/Fulcrum/CCVista/RO/<celex>>, unde <celex> este numărul unic de identificare al documentului). Numărul total de documente în limba română disponibile pe site-ul CCVista este de 19.286.

Fișierele au fost convertite din formatul Microsoft Word în formatul XML conform specificațiilor TEI. Conversia celor 19.286 fișiere a fost făcută automat, folosind pachetul de funcții „Visual Studio Tools for Office”. Acestea permit interacțiunea cu aplicația Microsoft Office direct din mediul de programare.

Datorită particularităților formatului, conversia documentelor a implicat și o serie de etape intermediare: au fost înlăturate comentariile traducătorilor, au fost șterse notele de subsol și secțiunile de cap de pagină, a fost normalizată folosirea caracterelor diacritice (unele documente foloseau „ș” și „ț” cu sedilă, iar altele „ș” și „ț” cu virgulă).

Dintre cele 19.286 de fișiere în format Microsoft Word au fost convertite 19.211 (restul de documente având erori de format). În formatul TEI-XML al documentelor românești au fost adăugate și datele de indexare Eurovoc acolo unde acestea erau disponibile. În (Ceaușu, 2008) este descris pe larg procesul de colectare și procesare a documentelor românești ale corpului JRC-Acquis.

3 Concluzii

Aplicațiile din lingvistica computațională, indiferent de tip, nu pot avea performanțe competitive în absența unor resurse lingvistice de foarte bună calitate. Aceste aplicații reprezintă, de cele mai multe ori, o combinație de module separate, fiecare cu propriile necesități informaționale. De exemplu, un sistem de întrebare-răspuns necesită o prelucrare primară a întrebării

utilizatorului pentru a putea extrage automat din aceasta o frază de interogare pentru motorul de căutare. La rândul ei, prelucrarea primară constă în adnotare morfolexicală și lematizare, procese care nu pot funcționa fără existența unor lexicoane sau modele de limbă, acestea din urmă extrase din corpusuri adnotate morfolexical. În consecință, resursele lingvistice computaționale reprezintă, în cele din urmă, o aproximare (materializată prin elaborarea de modele de limbă, lexicoane, ontologii lexicale etc.) a competențelor lingvistice ale vorbitorilor.

Calitatea resurselor lingvistice computaționale este direct responsabilă de performanțele algoritmilor de prelucrare automată a limbajului natural. Resursele prezentate în această lucrare au fost folosite de ICIA la competiția sistemelor de întrebare-răspuns pentru limba română CLEF 2009 în care Institutul nostru a ocupat locul întâi. Se confirmă faptul că aceste resurse lingvistice computaționale sunt valoroase și că orice efort de a le perfecționa reprezintă o foarte bună investiție din toate punctele de vedere.

Mulțumiri

Cercetările și rezultatele prezentate în această lucrare au fost susținute prin proiectul nr. D11-007 „SIR-RESDEC” din programul național de cercetare PN2.

Referințe

- Bentivogli, L., Forner, P., Magnini, B., Pianta, E. *Revising WordNet Domains Hierarchy: Semantics, Coverage, and Balancing*. În Proceedings of COLING 2004 Workshop on "Multilingual Linguistic Resources", Geneva, Switzerland, August 28, 2004, p. 101-108.
- Bloksma, L., Diez-Orzas, Vossen, P., *The User Requirements and Functional Specification of the EuroWordNet-project*. EWN-deliverable D.001, LE-4003, 1996.
- Ceașu, A. *Colectarea și procesarea documentelor românești ale corpusului JRC-Acquis*. În Resurse lingvistice și instrumente pentru prelucrarea limbii române, Iași, 2008, p. 125-130.
- Ceașu, A., Ștefănescu, D., Tufiș, D. *Acquis Communautaire sentence alignment using Support Vector Machines*. În Proceedings of the 5th LREC Conference, Genoa, Italy, 22-28 May, 2006.
- Chen, Q., Li, M., Zhou, M. *Improving Query Spelling Correction Using Web Search Results*. În Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning, Prague, 2007, p. 181-189.